

Generalization error bounds for stationary autoregressive models

Daniel J. McDonald
Department of Statistics
Carnegie Mellon University
Pittsburgh, PA 15213
danielmc@stat.cmu.edu

Cosma Rohilla Shalizi
Department of Statistics
Carnegie Mellon University
Pittsburgh, PA 15213
cshalizi@stat.cmu.edu

Mark Schervish
Department of Statistics
Carnegie Mellon University
Pittsburgh, PA 15213
mark@cmu.edu

Version: June 2, 2011

Abstract

We derive generalization error bounds for stationary univariate autoregressive (AR) models. We show that imposing stationarity is enough to control the Gaussian complexity without further regularization. This lets us use structural risk minimization for model selection. We demonstrate our methods by predicting interest rate movements.

1 Introduction

In standard machine learning situations, we observe one variable, X , and wish to predict another variable, Y , with an unknown joint distribution. Time series models are slightly different: we observe a sequence of observations $\mathbf{X}_1^n \equiv \{X_t\}_{t=1}^n$ from some process, and we wish to predict X_{n+h} , for some $h \in \mathbb{N}$. Throughout what follows, $\mathbf{X} = \{X_t\}_{t=-\infty}^{\infty}$ will be a sequence of random variables, i.e., each X_t is a measurable mapping from some probability space $(\Omega, \mathcal{F}, \mathbb{P})$ into a measurable space \mathcal{X} . A block of the random sequence will be written $\mathbf{X}_i^j \equiv \{X_t\}_{t=i}^j$, where either limit may go to infinity.

The goal in building a predictive model is to learn a function \hat{f} which maps the past into predictions for the future, evaluating the resulting forecasts through a loss function $\ell(X_{n+h}, \hat{f}(\mathbf{X}_1^n))$ which gives the cost of errors. Ideally, we would use f^* , the function which minimizes the risk

$$R(f) \equiv \mathbb{E}[\ell(X_{n+h}, f(\mathbf{X}_1^n))],$$

over all $f \in \mathcal{F}$, the class of prediction functions we can use.

Since the true joint distribution of the sequence is unknown, so is $R(f)$, but it is often estimated with the error on a training sample of size n

$$\hat{R}_n(f) \equiv \frac{1}{n} \sum_{t=1}^n \ell(X_{t+h}, f(\mathbf{X}_1^t)), \quad (1)$$

with \hat{f} being the minimizer of \hat{R}_n over \mathcal{F} . This is “empirical risk minimization”.

While $\hat{R}_n(\hat{f})$ converges to $R(\hat{f})$ for many algorithms, one can show that when \hat{f} minimizes (1), $\mathbb{E}[\hat{R}_n(\hat{f})] \leq R(\hat{f})$. This is because the choice of \hat{f} adapts to the training data, causing the

training error to be an over-optimistic estimate of the true risk. Also, training error must shrink as model complexity grows. Thus, empirical risk minimization gives unsatisfying results: it will tend to overfit the data and give poor out-of-sample predictions. Statistics and machine learning propose two mitigation strategies. The first is to restrict the class \mathcal{F} . The second, which we follow, is to change the optimization problem, penalizing model complexity. Without the true distribution, the prediction risk or generalization error are inaccessible. Instead, the goal is finding bounds on the risk which hold with high probability — “probably approximately correct” (PAC) bounds. A typical result is a confidence bound on the risk which says that with probability at least $1 - \eta$,

$$R(\hat{f}) \leq \hat{R}_n(\hat{f}) + \delta(C(\mathcal{F}), n, \eta),$$

where $C(\cdot)$ measures the complexity of the model class \mathcal{F} , and $\delta(\cdot)$ is a function of this complexity, the confidence level, and the number of observed data points.

The statistics and machine learning literature contains many generalization error bounds for both classification and regression problems with IID data, but their extension to time series prediction is a fairly recent development; in 1997, Vidyasagar [21] named extending such results to time series as an important open problem. Yu [22] sets forth many of the uniform ergodic theorems that are needed to derive generalization error bounds for stochastic processes. Meir [11] is one of the first papers to construct risk bounds for time series. His approach was to consider a stationary but infinite-memory process, and to decompose the training error of a predictor with finite memory, chosen through empirical risk minimization, into three parts:

$$\hat{R}(\hat{f}_{p,n,d}) = (\hat{R}(\hat{f}_{p,n,d}) - \hat{R}(f_{p,n}^*)) + (\hat{R}(f_{p,n}^*) - \hat{R}(f_p^*)) + \hat{R}(f_p^*)$$

where $\hat{f}_{p,n,d}$ is an empirical estimate based on finite data of length n , finite memory of length p , and complexity indexed by d ; $f_{p,d}^*$ is the oracle with finite memory and given complexity, and f_p^* is the oracle with finite memory over all possible complexities. The three terms amount to an estimation error incurred from the use of limited and noisy data, an approximation error due to selecting a predictor from a class of limited complexity, and a loss from approximating an infinite memory process with a finite memory process.

More recently, others have provided PAC results for non-IID data. Steinwart and Christmann [20] prove an oracle inequality for generic regularized empirical risk minimization algorithms learning from α -mixing processes, a fairly general sort of weak serial dependence, getting learning rates for least-squares support vector machines (SVMs) close to the optimal IID rates. Mohri and Ros-tamizadeh [13] prove stability-based generalization bounds when the data are stationary and φ -mixing or β -mixing, strictly generalizing IID results and applying to all stable learning algorithms. (We define β -mixing below.) Karandikar and Vidyasagar [8] show that if an algorithm is “sub-additive” and yields a predictor whose risk can be upper bounded when the data are IID, then the same algorithm yields predictors whose risk can be bounded if data are β -mixing. They use this result to derive generalization error bounds in terms of the learning rates for IID data and the β -mixing coefficients.

All these generalization bounds for dependent data rely on notions of complexity which, while common in machine learning, are hard to apply to models and algorithms ubiquitous in the time series literature. SVMs, neural networks, and kernel methods have known complexities, so their risk can be bounded on dependent data as well. On the other hand, autoregressive moving average (ARMA) models, generalized autoregressive conditional heteroskedasticity (GARCH) models, and state-space models in general have unknown complexity and are therefore neglected theoretically. (This does not keep them from being used in applied statistics, or even in machine learning and robotics, e.g., [17, 15, 18, 2, 10].) Arbitrarily regularizing such models will not do, as often the only assumption applied researchers are willing to make is that the time series is stationary.

We show that the assumption of stationarity regularizes autoregressive (AR) models implicitly, allowing for the application of risk bounds without the need for additional penalties. This result follows from work in the optimal control and systems design literatures but the application is novel. In §2, we introduce concepts from time series and complexity theory necessary for our results. Section 3 uses these results to calculate explicit risk bounds for autoregressive models. Section 4 illustrates

the applicability of our methods by forecasting interest rate movements. We discuss our results and articulate directions for future research in §5.

2 Preliminaries

Before developing our results, we need to explain the idea of effective sample size for dependent data, and the closely related measure of serial dependence called β -mixing, as well as the Gaussian complexity technique for measuring model complexity.

2.1 Time series

Because time-series data are dependent, the number of data points n in a sample \mathbf{X}_1^n exaggerates how much information the sample contains. Knowing the past allows forecasters to predict future data (at least to some degree), so actually observing those future data points gives less information about the underlying data generating process than in the IID case. Thus, the sample size term in a probabilistic risk bound must be adjusted to reflect the dependence in the data source. This effective sample size may be much less than n .

We investigate only stationary β -mixing input data. We first remind the reader of the notion of (strict or strong) stationarity.

Definition 2.1 (Stationarity). *A sequence of random variables \mathbf{X} is stationary when all its finite-dimensional distributions are invariant over time: for all t and all non-negative integers i and j , the random vectors \mathbf{X}_t^{t+i} and \mathbf{X}_{t+j}^{t+i+j} have the same distribution.*

From among all the stationary processes, we restrict ourselves to ones where widely-separated observations are asymptotically independent. Stationarity does not imply that the random variables X_t are independent across time t , only that the distribution of X_t is constant in time. The next definition describes the nature of the serial dependence which we are willing to allow.

Definition 2.2 (β -Mixing). *Let $\sigma_i^j = \sigma(\mathbf{X}_i^j)$ be the σ -field of events generated by the appropriate collection of random variables. Let \mathbb{P}_t be the restriction of \mathbb{P} to $\sigma_{-\infty}^t$, \mathbb{P}_{t+m} be the restriction of \mathbb{P} to σ_{t+m}^∞ , and $\mathbb{P}_{t \otimes t+m}$ be the restriction of \mathbb{P} to $\sigma(\mathbf{X}_{-\infty}^t, \mathbf{X}_{t+m}^\infty)$. The coefficient of absolute regularity, or β -mixing coefficient, $\beta(m)$, is given by*

$$\beta(m) \equiv \|\mathbb{P}_t \times \mathbb{P}_{t+m} - \mathbb{P}_{t \otimes t+m}\|_{TV}, \quad (2)$$

where $\|\cdot\|_{TV}$ is the total variation norm. A stochastic process is absolutely regular, or β -mixing, if $\beta(m) \rightarrow 0$ as $m \rightarrow \infty$.

This is only one of many equivalent characterizations of β -mixing (see Bradley [3] for others). This definition makes clear that a process is β -mixing if the joint probability of events which are widely separated in time increasingly approaches the product of the individual probabilities, i.e., that \mathbf{X} is asymptotically independent. Typically, a supremum over t is taken in (2), however, this is unnecessary for stationary processes, i.e. $\beta(m)$ as defined above is independent of t .

2.2 Gaussian complexity

Statistical learning theory provides several ways of measuring the complexity of a class of predictive models. The results we are using here rely on Gaussian complexity (see, e.g., Bartlett and Mendelson [1]), which can be thought of as measuring how well the model can (seem to) fit white noise.

Definition 2.3 (Gaussian Complexity). *Let \mathbf{X}_1^n be a (not necessarily IID) sample drawn according to ν . The empirical Gaussian complexity is*

$$\widehat{\mathfrak{G}}_n(\mathcal{F}) \equiv 2\mathbb{E}_Z \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n Z_i f(\mathbf{X}_1^i) \right| \mid \mathbf{X}_1^n \right],$$

where Z_i are a sequence of random variables, independent of each other and everything else, and drawn from a standard Gaussian distribution. The Gaussian complexity is

$$\mathfrak{G}_n(\mathcal{F}) \equiv \mathbb{E}_\nu \left[\widehat{\mathfrak{G}}_n(\mathcal{F}) \right]$$

where the expectation is over sample paths D_n generated by ν .

The term inside the supremum, $|\frac{1}{n} \sum_{i=1}^n Z_i f(\mathbf{X}_1^i)|$, is the sample covariance between the noise Z and the predictions of a particular model f . The Gaussian complexity takes the largest value of this sample covariance over all models in the class (mimicking empirical risk minimization), then averages over realizations of the noise.

Intuitively, Gaussian complexity measures how well our models could seem to fit outcomes which were really just noise, giving a baseline against which to assess the risk of over-fitting or failing to generalize. As the sample size n grows, for any given f the sample covariance $|\frac{1}{n} \sum_{i=1}^n Z_i f(\mathbf{X}_1^i)| \rightarrow 0$, by the ergodic theorem; the overall Gaussian complexity should also shrink, though more slowly, unless the model class is so flexible that it can fit absolutely anything, in which case one can conclude nothing about how well it will predict in the future from the fact that it performed well in the past.

2.3 Error bounds for β -mixing data

Mohri and Rostamizadeh [12] present Gaussian¹ complexity-based error bounds for stationary β -mixing sequences, a generalization of similar bounds presented earlier for the IID case. The results are data-dependent and measure the complexity of a class of hypotheses based on the training sample.

Theorem 2.4. *Let \mathcal{F} be a space of candidate predictors and let \mathcal{H} be the space of induced losses:*

$$\mathcal{H} = \{h = \ell(\cdot, f(\cdot)) : f \in \mathcal{F}\}$$

for some loss function $0 \leq \ell(\cdot, \cdot) \leq M$. Then for any sample \mathbf{X}_1^n drawn from a stationary β -mixing distribution, and for any $\mu, m > 0$ with $2\mu m = n$ and $\eta > 4(\mu - 1)\beta(m)$ where $\beta(m)$ is the mixing coefficient, with probability at least $1 - \eta$,

$$R(\widehat{f}) \leq \widehat{R}_n(\widehat{f}) + \left(\frac{\pi}{2}\right)^{1/2} \widehat{\mathfrak{G}}_\mu(\mathcal{H}) + 3M \sqrt{\frac{\ln 4/\eta'}{2\mu}},$$

and

$$R(\widehat{f}) \leq \widehat{R}_n(\widehat{f}) + \left(\frac{\pi}{2}\right)^{1/2} \mathfrak{G}_\mu(\mathcal{H}) + M \sqrt{\frac{\ln 2/\eta'}{2\mu}},$$

where $\eta' = \eta - 4(\mu - 1)\beta(m)$ in the first case or $\eta' = \eta - 2(\mu - 1)\beta(m)$ in the second. .

The generalization error bounds in Theorem 2.4 have a straightforward interpretation. The risk of a chosen model is controlled, with high probability, by three terms. The first term, the training error, describes how well the model performs in-sample. More complicated models can more closely fit any data set, so increased complexity leads to smaller training error. This is penalized by the second term, the Gaussian complexity. The first bound uses the empirical Gaussian complexity which is calculated from the data \mathbf{X}_1^n while the second uses the expected Gaussian complexity, and is therefore tighter. The third term is the confidence term and is a function only of the confidence level η and the effective number of data points on which the model was based μ . While it was actually trained on n data points, because of dependence, this number must be reduced. This process is accomplished by taking μ widely spaced blocks of points. Under the asymptotic independence quantified by β , this spacing lets us treat these blocks as independent.

¹In fact, they present the bounds in terms of the Rademacher complexity, a closely related idea. However, using Gaussian complexity instead requires no modifications to their results while simplifying the proofs contained here. The constant $(\pi/2)^{1/2}$ in Theorem 2.4 is given in Ledoux and Talagrand [9].

3 Results

Autoregressive models are used frequently in economics, finance, and other disciplines. Their main utility lies in their straightforward parametric form, as well as their interpretability: predictions for the future are linear combinations of some fixed length of previous observations. See Shumway and Stoffer [19] for a standard introduction.

Suppose that \mathbf{X} is a real-valued random sequence, evolving as

$$X_t = \sum_{i=1}^p \phi_i X_{t-i} + \epsilon_t,$$

where ϵ_t has mean zero, finite variance, $\epsilon_j \perp \epsilon_i$ for all $i \neq j$, and $\epsilon_i \perp X_j$ for all $i > j$. This is the traditional specification of an *autoregressive order p* or $\text{AR}(p)$ model. Having observed data $\{X_t\}_{t=1}^n$, and supposing p to be known, fitting the model amounts to estimating the coefficients $\{\phi\}_{i=1}^p$. The most natural way to do this is to use ordinary least squares (OLS). Let

$$\phi = \begin{pmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_p \end{pmatrix} \quad \mathbb{Y} = \begin{pmatrix} X_{p+1} \\ X_{p+2} \\ \vdots \\ X_{n-1} \\ X_n \end{pmatrix} \quad \mathbb{X} = \begin{pmatrix} X_p & X_{p-1} & \cdots & X_1 \\ X_{p+1} & X_p & \cdots & X_2 \\ \vdots & \vdots & \ddots & \vdots \\ X_{n-2} & X_{n-3} & \cdots & X_{n-p-1} \\ X_{n-1} & X_{n-2} & \cdots & X_{n-p} \end{pmatrix}.$$

Generalization error bounds for these processes follow from an ability to characterize their Gaussian complexity. The theorem below uses stationarity to bound the risk of AR models. The remainder of this section provides the components necessary to prove the results.

Theorem 3.1. *Let D_n be a sample of length n from a stationary β -mixing distribution. For any $\mu, m > 0$ with $2\mu m = n$ and $\eta > 4(\mu - 1)\beta(m)$, then under squared error loss truncated at M , the prediction error of an $\text{AR}(p)$ ($p > 1$) model can be bounded with probability at least $1 - \eta$ using*

$$R(\hat{f}) \leq \hat{R}_n(\hat{f}) + \frac{4\sqrt{\pi M \log(p+1)}}{\mu} \max_{1 \leq j, j' \leq p+1} \left(\sum_{i \in \mathcal{I}} \langle \mathbb{X}_i, \phi_j - \phi_{j'} \rangle^2 \right)^{1/2} + 3M \sqrt{\frac{\ln 4/\eta'}{2\mu}},$$

or

$$R(\hat{f}) \leq \hat{R}_n(\hat{f}) + \frac{4\sqrt{\pi M \log(p+1)}}{\mu} \mathbb{E} \left[\max_{1 \leq j, j' \leq p+1} \left(\sum_{i \in \mathcal{I}} \langle \mathbb{X}_i, \phi_j - \phi_{j'} \rangle^2 \right)^{1/2} \right] + M \sqrt{\frac{\ln 2/\eta'}{2\mu}},$$

where $\mathcal{I} = \{i : i = \lfloor a/2 \rfloor + 2ak, 0 \leq k \leq \mu\}$, ϕ_j is the j^{th} vertex of the stability domain, and \mathbb{X}_i is the i^{th} row of the design matrix \mathbb{X} .

For $p = 1$ slight adjustments are required. We state this result as a corollary.

Corollary 3.2. *Under the same conditions as above, the prediction error of an $\text{AR}(1)$ model can be bounded with probability at least $1 - \eta$ using*

$$R(\hat{f}) \leq \hat{R}_n(\hat{f}) + \frac{4}{\mu} \sqrt{\frac{M}{2}} \left(\sum_{i \in \mathcal{I}} X_i^2 \right)^{1/2} + 3M \sqrt{\frac{\ln 4/\eta'}{2\mu}},$$

or

$$R(\hat{f}) \leq \hat{R}_n(\hat{f}) + \frac{4}{\mu} \sqrt{\frac{M}{2}} \mathbb{E} \left[\left(\sum_{i \in \mathcal{I}} \mathbb{X}_i^2 \right)^{1/2} \right] + M \sqrt{\frac{\ln 2/\eta'}{2\mu}}.$$

3.1 Proof components

To prove Theorem 3.1 it is necessary to control the size of the model class by using the stationarity assumption.

3.1.1 Stationarity controls the hypothesis space

Define, as an estimator of ϕ ,

$$\hat{\phi} \equiv \underset{\phi}{\operatorname{argmin}} \|\mathbb{Y} - \mathbb{X}\phi\|_2^2, \quad (3)$$

where $\|\cdot\|_2$ is the Euclidean norm.² Equation 3 has the usual closed form OLS solution:

$$\hat{\phi} = (\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\mathbb{Y}. \quad (4)$$

Despite the simplicity of Eq. 4, modellers often require that the estimated autoregressive process be stationary. This can be checked algebraically: the complex roots of the polynomial

$$Q_p(z) = z^p - \phi_1 z^{p-1} - \phi_2 z^{p-2} - \dots - \phi_p$$

must lie strictly inside the unit circle. Eq. 3 is thus not quite right for estimating a stationary autoregressive model, as it does not incorporate this constraint.

Constraining the roots of $Q_p(z)$ constrains the coefficients ϕ . The set ϕ where the process is stationary is the stability domain, \mathcal{B}_p . Clearly, \mathcal{B}_1 is just $|\phi_1| < 1$. Fam and Meditch [6] gives a recursive method for determining \mathcal{B}_p for general p . In particular, they show that the convex hull of the space of stationary solutions is a convex polyhedron with vertices at the extremes of the \mathcal{B}_p . This convex hull basically determines the complexity of stationary AR models.

3.1.2 Gaussian complexity of AR models

Returning to the $\text{AR}(p)$ model, it is necessary to find the Gaussian complexity of the function class

$$\mathcal{F}_p = \left\{ \phi : x_t = \sum_{i=1}^p \phi_i x_{t-i} \text{ and } x_t \text{ is stationary} \right\}.$$

Theorem 3.3. *For the $\text{AR}(p)$ model with $p > 1$, the empirical Gaussian complexity is given by*

$$\hat{\mathfrak{G}}_k(\mathcal{F}) \leq \frac{2\sqrt{2}}{n} (\log(p+1))^{1/2} \max_{1 \leq j, j' \leq p+1} \left(\sum_{i=1}^k \langle \mathbb{X}_i, \phi_j - \phi_{j'} \rangle^2 \right)^{1/2},$$

where ϕ_j is the j^{th} vertex of the stability domain and \mathbb{X}_i is the i^{th} row of the design matrix \mathbb{X} .

The proof relies on the following version of Slepian's Lemma (see, for example Ledoux and Talagrand [9] or Bartlett and Mendelson [1]).

Lemma 3.4 (Slepian). *Let V_1, \dots, V_k be random variables such that for all $1 \leq j \leq k$, $V_j = \sum_{i=1}^n a_{ij} g_i$ where g_1, \dots, g_n are iid standard normal random variables. Then,*

$$\mathbb{E}[\max_j V_j] \leq \sqrt{2}(\log k)^{1/2} \max_{j, j'} \sqrt{\mathbb{E}[(V_j - V_{j'})^2]}.$$

Proof of Theorem 3.3.

$$\begin{aligned} \hat{\mathfrak{G}}_n \mathcal{F} &= \mathbb{E} \sup_{f \in \mathcal{F}} \frac{2}{n} \sum_{i=1}^n g_i f(x_i) = \mathbb{E} \sup_{\phi \in \mathcal{B}_p} \frac{2}{n} \sum_{i=1}^n g_i \langle \mathbb{X}_i, \phi \rangle \\ &= \mathbb{E} \sup_{\phi \in \mathcal{B}_p} \left\langle \frac{2}{n} \sum_{i=1}^n g_i \mathbb{X}_i, \phi \right\rangle = \mathbb{E} \sup_{\phi \in \operatorname{conv}(\mathcal{B}_p)} \left\langle \frac{2}{n} \sum_{i=1}^n g_i \mathbb{X}_i, \phi \right\rangle, \end{aligned}$$

²There are other ways to estimate AR models, but they typically amount to very similar optimization problems.

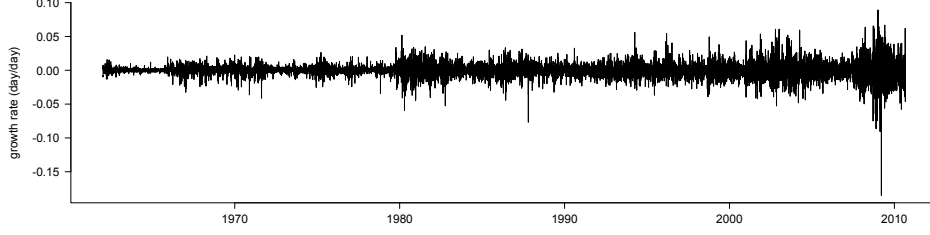


Figure 1: Growth rate of 10-year treasury bond

where the last equality follows from Theorem 12 in [1]. By standard results from convex optimization, this supremum is attained at one of the vertices of $\text{conv}(\mathcal{B}_p)$. Therefore,

$$\hat{\mathfrak{G}}_n(\mathcal{F}) = \mathbb{E}\left[\max_j \frac{2}{n} \sum_{i=1}^n g_i \langle \mathbb{X}_i, \phi_j \rangle\right],$$

where ϕ_j is the j^{th} vertex of $\text{conv}(\mathcal{B}_p)$. Let $V_j = \sum_{i=1}^n g_i \langle \mathbb{X}_i, \phi_j \rangle$. Then by the Lemma 3.4,

$$\begin{aligned} \hat{\mathfrak{G}}_n(\mathcal{F}) &\leq \frac{2\sqrt{2}}{n} (\log p + 1)^{1/2} \max_{j,j'} \sqrt{\mathbb{E}[V_j - V_{j'}]^2} \\ &= \frac{2\sqrt{2}}{n} (\log p + 1)^{1/2} \max_{j,j'} \sqrt{\mathbb{E}\left[\sum_{i=1}^n g_i \langle \mathbb{X}_i, \phi_j - \phi_{j'} \rangle\right]^2} \\ &= \frac{2\sqrt{2}}{n} (\log p + 1)^{1/2} \max_{1 \leq j, j' \leq p} \sqrt{\sum_{i=1}^n \langle \mathbb{X}_i, \phi_j - \phi_{j'} \rangle^2} \end{aligned}$$

where \mathbb{X}_i is the i^{th} -entry of the design matrix. \square

When $p = 1$, as in Corollary 3.2, we can calculate the complexity directly. The proof's last line shows that we are essentially interested in the diameter of the stability domain \mathcal{B}_p projected onto the column space of \mathbb{X} , which gives a tighter bound than that from the general results on linear prediction in e.g. Kakade et al. [7].

Since we care about the complexity of the model class \mathcal{F} viewed through the loss function ℓ , we must also account for this additional complexity. For c -Lipschitz loss functions, this just means multiplying $\mathfrak{G}_n(\mathcal{F})$ by $2c$.

4 Application

We illustrate our results by predicting interest rate changes — specifically, the 10-year Treasury Constant Maturity Rate series from the Federal Reserve Bank of St. Louis' FRED database³ — recorded daily from January 2, 1962 to August 31, 2010. Transforming the series into daily natural-log growth rates leaves $n = 12150$ observations (Figure 1). The changing variance apparent in the figure is why interest rates are typically forecast with GARCH(1,1) models. For this illustration however, we will use an AR(p) model, picking the memory order p by the risk bound.

Figure 2 shows the training error

$$\hat{R}_n(\hat{f}) = \frac{1}{n-p} \sum_{t=p+1}^n (\hat{X}_t - X_t)^2$$

³Available at <http://research.stlouisfed.org/fred2/series/DGS10?cid=115>.

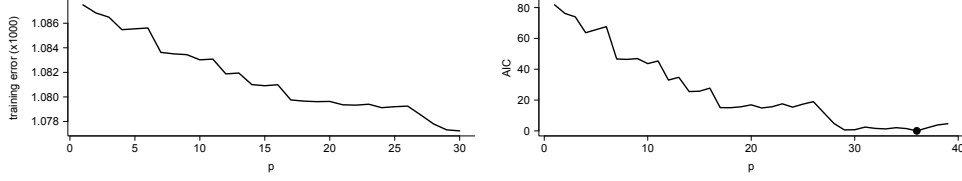


Figure 2: Training error (top panel) and AIC (bottom panel) against model order

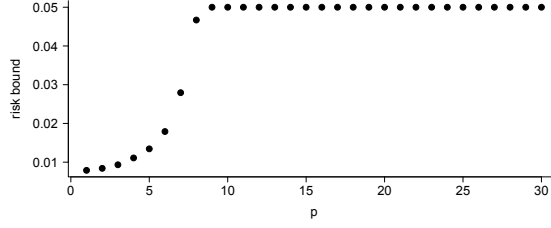


Figure 3: Generalization error bound for different model orders

where X_t is the t^{th} datapoint, and \hat{X}_t is the model's prediction. \hat{R}_n shrinks as the order of the model (p) grows, as it must since ordinary least squares minimizes \hat{R}_n for a given p . Also shown is the gap between the AIC for different p and the lowest attainable value; this would select an AR(36) model.

A better strategy uses the probabilistic risk bound derived above. The goal of model selection is to pick, with high probability, the model with the smallest risk; this is Vapnik's structural risk minimization principle. Here, it is clear that AIC is dramatically overfitting. The optimal model using the risk bound is an AR(1). Figure 3 plots the risk bound against p with the loss function truncated at 0.05. (No daily interest rate change has ever had loss larger than 0.034, and results are fairly insensitive to the level of the loss cap.) This bound says that with 95% probability, *regardless of the true data generating process*, the AR(1) model will make mistakes with squared error no larger than 0.0079. If we had instead predicted with zero, this loss would have occurred three times.

One issue with Theorem 2.4 is that it requires knowledge of the β -mixing coefficients, $\beta(m)$. Of course, the dependence structure of this data is unknown, so we calculated it under generous assumptions on the data generating process. In a homogeneous Markov process, the β -mixing coefficients work out to

$$\beta(m) = \int \pi(dx) \|P^m(x, \cdot) - \pi\|_{TV}$$

where $P^m(x, \cdot)$ is the m -step transition operator and π is the stationary distribution [14, 4]. Since AR models are Markovian, we estimated an AR(q) model with Gaussian errors for q large and calculated the mixing coefficients using the stationary and transition distributions. To create the bound, we used $m = 7$ and $\mu = 867$. We address non-parametric estimation of β -mixing coefficients elsewhere [Anon.].

5 Discussion

We have constructed a finite-sample predictive risk bound for autoregressive models, using the stationarity assumption to constrain OLS estimation. Interestingly, stationarity — a common assumption among applied researchers — constrains the model space enough to yield bounds without further regularization. Moreover, this is the first predictive risk bound we know of for any of the standard models of time series analysis.

Traditionally, time series analysts have selected models by blending empirical risk minimization, more-or-less quantitative inspection of the residuals (e.g., the Box-Ljung test; see [19]), and AIC. In many applications, however, what really matters is prediction, and none of these techniques, including AIC, controls generalization error, especially with mis-specification. (Cross-validation is a partial exception, but it is tricky for time series; see [16] and references therein.) Our bound controls prediction risk directly. Admittedly, our bound covers only univariate autoregressive models, the plainest of a large family of traditional time series models, but we believe a similar result will cover the more elaborate members of the family such as vector autoregressive, autoregressive-moving average, or autoregressive conditionally heteroskedastic models. While the characterization of the stationary domain from [6] on which we relied breaks down for such models, they are all variants of the linear state space model [5], whose parameters are restricted under stationarity, and so we hope to obtain a general risk bound, possibly with stronger variants for particular specifications.

References

- [1] Peter L. Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- [2] B.C. Becker, H. Tummala, and C.N. Riviere. Autoregressive modeling of physiological tremor under microsurgical conditions. In *Engineering in Medicine and Biology Society, 2008. EMBS 2008. 30th Annual International Conference of the IEEE*, pages 1948–1951. IEEE, 2008.
- [3] Richard C. Bradley. Basic properties of strong mixing conditions. a survey and some open questions. *Probability Surveys*, 2:107–144, 2005. URL <http://arxiv.org/abs/math/0511078>.
- [4] Y.A. Davydov. Mixing conditions for markov chains. *Theory of Probability and its Applications*, 18(2):312–328, 1973.
- [5] J. Durbin and S.J. Koopman. *Time Series Analysis by State Space Methods*. Oxford Univ Press, Oxford, 2001.
- [6] Adly T. Fam and James S. Meditch. A canonical parameter space for linear systems design. *IEEE Transactions on Automatic Control*, 23(3):454–458, 1978.
- [7] Sham M. Kakade, Karthik Sridharan, and Ambuj Tewari. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. Technical report, NIPS, 2008. URL <http://ttic.uchicago.edu/~karthik/rad-paper.pdf>.
- [8] R. L. Karandikar and M. Vidyasagar. Probably approximately correct learning with beta-mixing input sequences. submitted for publication, 2009.
- [9] M. Ledoux and M. Talagrand. *Probability in Banach Spaces: Isoperimetry and Processes*. A Series of Modern Surveys in Mathematics. Springer Verlag, Berlin, 1991. ISBN 3540520139.
- [10] J. Li and A.W. Moore. Forecasting web page views: Methods and observations. *Journal of Machine Learning Research*, 9:2217–2250, 2008.
- [11] Ron Meir. Nonparametric time series prediction through adaptive model selection. *Machine Learning*, 39(1):5–34, 2000. URL <http://www.ee.technion.ac.il/~rmeir/Publications/MeirTimeSeries00.pdf>.
- [12] Mehryar Mohri and Afshin Rostamizadeh. Rademacher complexity bounds for non-iid processes. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *Advances in Neural Information Processing Systems 21*, volume 21, pages 1097–1104, 2009.
- [13] Mehryar Mohri and Afshin Rostamizadeh. Stability bounds for stationary φ -mixing and β -mixing processes. *Journal of Machine Learning Research*, 11:789–814, February 2010.

- [14] A. Mokkadem. Mixing properties of arma processes. *Stochastic processes and their applications*, 29(2):309–315, 1988.
- [15] R.K. Olsson and L.K. Hansen. Linear state-space models for blind source separation. *The Journal of Machine Learning Research*, 7:2585–2602, 2006. ISSN 1532-4435.
- [16] J. Racine. Consistent cross-validatory model-selection for dependent data: Hv-block cross-validation. *Journal of econometrics*, 99(1):39–61, 2000.
- [17] J. Ruiz-del Solar and P. Vallejos. Motion detection and tracking for an aibo robot using motion compensation and kalman filtering. In *Lecture Notes in Computer Science 3276 (RoboCup 2004)*, pages 619–627. Springer Verlag, 2005.
- [18] M. Sak, D.L. Dowe, and S. Ray. Minimum message length moving average time series data mining. In *Computational Intelligence Methods and Applications, 2005 ICSC Congress on*, page 6. IEEE, 2006. ISBN 1424400201.
- [19] R.H. Shumway and D.S. Stoffer. *Time Series Analysis and Its Applications*. Springer Series in Statistics. Springer Verlag, New York, 2000.
- [20] Ingo Steinwart and Andreas Christmann. Fast learning from non-i.i.d. observations. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1768–1776. MIT Press, 2009. URL http://books.nips.cc/papers/files/nips22/NIPS2009_1061.pdf.
- [21] M. Vidyasagar. *A Theory of Learning and Generalization: With Applications to Neural Networks and Control Systems*. Springer Verlag, Berlin, 1997.
- [22] Bin Yu. Rates of convergence for empirical processes of stationary mixing sequences. *The Annals of Probability*, 22(1):94–116, 1994.